



UNITED STATES PATENT AND TRADEMARK OFFICE

UNITED STATES DEPARTMENT OF COMMERCE
United States Patent and Trademark Office
Address: COMMISSIONER FOR PATENTS
P.O. Box 1450
Alexandria, Virginia 22313-1450
www.uspto.gov

APPLICATION NO.	FILING DATE	FIRST NAMED INVENTOR	ATTORNEY DOCKET NO.	CONFIRMATION NO.
09/978,192	10/15/2001	Avi J. Ashkenazi	GNE.2630PIC9	3437

7590 12/15/2004

Ginger R. Dreger, Esq.
Knobbe Martens Olson & Bear
620 NEWPORT CENTER DRIVE
SIXTEENTH FLOOR
NEWPORT BEACH, CA 92660

EXAMINER

O HARA, EILEEN B

ART UNIT	PAPER NUMBER
----------	--------------

1646

DATE MAILED: 12/15/2004

Please find below and/or attached an Office communication concerning this application or proceeding.

Office Action Summary

Application No.

09/978,192

Applicant(s)

ASHKENAZI ET AL.

Examiner

Eileen O'Hara

Art Unit

1646

-- The MAILING DATE of this communication appears on the cover sheet with the correspondence address --

Period for Reply

A SHORTENED STATUTORY PERIOD FOR REPLY IS SET TO EXPIRE 3 MONTH(S) FROM THE MAILING DATE OF THIS COMMUNICATION.

- Extensions of time may be available under the provisions of 37 CFR 1.136(a). In no event, however, may a reply be timely filed after SIX (6) MONTHS from the mailing date of this communication.
- If the period for reply specified above is less than thirty (30) days, a reply within the statutory minimum of thirty (30) days will be considered timely.
- If NO period for reply is specified above, the maximum statutory period will apply and will expire SIX (6) MONTHS from the mailing date of this communication.
- Failure to reply within the set or extended period for reply will, by statute, cause the application to become ABANDONED (35 U.S.C. § 133). Any reply received by the Office later than three months after the mailing date of this communication, even if timely filed, may reduce any earned patent term adjustment. See 37 CFR 1.704(b).

Status

- 1) ☒ Responsive to communication(s) filed on 14 September 2004.
- 2a) ☒ This action is **FINAL**. 2b) ☐ This action is non-final.
- 3) ☐ Since this application is in condition for allowance except for formal matters, prosecution as to the merits is closed in accordance with the practice under *Ex parte Quayle*, 1935 C.D. 11, 453 O.G. 213.

Disposition of Claims

- 4) ☒ Claim(s) 58-62 is/are pending in the application.
- 4a) Of the above claim(s) _____ is/are withdrawn from consideration.
- 5) ☐ Claim(s) _____ is/are allowed.
- 6) ☒ Claim(s) 58-62 is/are rejected.
- 7) ☐ Claim(s) _____ is/are objected to.
- 8) ☐ Claim(s) _____ are subject to restriction and/or election requirement.

Application Papers

- 9) ☐ The specification is objected to by the Examiner.
- 10) ☒ The drawing(s) filed on 15 October 2001 is/are: a) ☒ accepted or b) ☐ objected to by the Examiner.
Applicant may not request that any objection to the drawing(s) be held in abeyance. See 37 CFR 1.85(a).
Replacement drawing sheet(s) including the correction is required if the drawing(s) is objected to. See 37 CFR 1.121(d).
- 11) ☐ The oath or declaration is objected to by the Examiner. Note the attached Office Action or form PTO-152.

Priority under 35 U.S.C. § 119

- 12) ☐ Acknowledgment is made of a claim for foreign priority under 35 U.S.C. § 119(a)-(d) or (f).
- a) ☐ All b) ☐ Some * c) ☐ None of:
- ☐ Certified copies of the priority documents have been received.
 - ☐ Certified copies of the priority documents have been received in Application No. _____.
 - ☐ Copies of the certified copies of the priority documents have been received in this National Stage application from the International Bureau (PCT Rule 17.2(a)).
- * See the attached detailed Office action for a list of the certified copies not received.

Attachment(s)

- | | |
|--|---|
| 1) <input type="checkbox"/> Notice of References Cited (PTO-892) | 4) <input type="checkbox"/> Interview Summary (PTO-413)
Paper No(s)/Mail Date. _____ |
| 2) <input type="checkbox"/> Notice of Draftsperson's Patent Drawing Review (PTO-948) | 5) <input type="checkbox"/> Notice of Informal Patent Application (PTO-152) |
| 3) <input type="checkbox"/> Information Disclosure Statement(s) (PTO-1449 or PTO/SB/08)
Paper No(s)/Mail Date _____ | 6) <input type="checkbox"/> Other: _____ |

DETAILED ACTION

1. Claims 58-62 are pending in the instant application. Claim 63 has been canceled and claim 58 has been amended as requested by Applicant in the amendment filed September 14, 2004.

Withdrawn Objections and Rejections

2. Any objection or rejection of record which is not expressly repeated in this action has been overcome by Applicant's response and withdrawn.

Maintained Rejections

Claim Rejections - 35 USC § 101 and § 112

35 U.S.C. 101 reads as follows:

Whoever invents or discovers any new and useful process, machine, manufacture, or composition of matter, or any new and useful improvement thereof, may obtain a patent therefor, subject to the conditions and requirements of this title.

3. Claims 58-62 remain rejected under 35 U.S.C. 101 because the claimed invention is not supported by either a specific and substantial asserted utility or a well established utility, for reasons of record in the previous office action, mailed May 20, 2004, at pages 4-8 and below.

Applicants' arguments (pages 10-19, Paper filed Sept. 14, 2004) have been fully considered but are not deemed persuasive.

Applicants traverse the rejection and discuss the legal standard for utility on pages 10-11, and starting on page 12 discuss the proper application of the legal standard. Applicants rely on the gene amplification data for patentable utility for the PRO274 protein and antibodies thereof,

Art Unit: 1646

and explain the gene amplification assay of Example 114, in which PRO274 is amplified more than two fold in three types of human primary lung tumors, which Applicants assert is significant and that the PRO274 gene has utility as a diagnostic of lung cancer. Applicants provide the Declaration by Dr. Audrey Goddard, in which she states that a gene identified as being amplified at least 2-fold by the quantitative TaqMan PCR assay in a tumor sample relative to a normal sample is useful as a marker for the diagnosis of cancer. Applicants assert that as the TaqMan realtime PCR method has gained wide recognition for its versatility, sensitivity and accuracy, and is in extensive use for the study of gene amplification, one of ordinary skill in the art would find it credible that PRO274 is a diagnostic marker of human lung cancer.

The Goddard Declaration filed under 37 CFR 1.132, filed Sept. 14, 2004 is insufficient to overcome the rejection of claims 58-62 as set forth in the last Office action because: while the declaration and supporting references are convincing that the TaqMan realtime PCR method is very sensitive and can identify amplified genes, the claims are drawn to antibodies to protein encoded by the PRO247 gene, and as discussed in the previous office action and below, it is not predictable that gene amplification results in increased mRNA expression, or that increased mRNA expression results in increased protein production.

Applicant argues that the Gygi et al. publication does not support the rejection. Applicant characterizes Gygi et al. as teaching that there is a general trend but no strong correlation between polypeptide expression level and transcript level. Applicant further characterizes Gygi et al.'s conclusions as showing that there is a positive correlation between transcript and polypeptide for most of the 150 yeast polypeptides studied, but the correlation is not linear and thus one cannot accurately predict polypeptide levels from mRNA levels. Applicant concludes

Art Unit: 1646

that Gygi et al. show that it is more likely than not that a positive correlation exists between mRNA and polypeptide levels. This has been fully considered but is not found to be persuasive. In the instant case, the specification provides data showing a very small increase in DNA copy number, approximately **2-fold**, in a few tumor samples for PRO274. There is no evidence regarding whether or not the PRO274 **mRNA** or **polypeptide** levels are also increased in these tumor samples. Since the instant claims are directed to **antibodies** to PRO274 **polypeptide**, it was imperative to find evidence in the relevant scientific literature whether or not a small increase in DNA copy number would be considered by the skilled artisan to be predictive of increased mRNA and polypeptide levels. Pennica et al. was cited as evidence showing a lack of correlation between gene (DNA) amplification and elevated mRNA levels. Gygi et al. was cited as providing evidence that polypeptide levels cannot be accurately predicted from mRNA levels, and that variances as much as **40-fold** or even **50-fold** were not uncommon. Given the small magnitude by which the DNA copy number of PRO274 is increased, and the evidence provided by Gygi et al. and Pennica et al., it is clear that one skilled in the art would not assume that a small increase in gene copy number would correlate with significantly increased mRNA or polypeptide levels. One skilled in the art would do further research to determine whether or not the PRO274 polypeptide levels increased significantly in the tumor samples. The requirement for such further research requirements makes it clear that the asserted utility is not yet in currently available form, i.e., it is not substantial. This further experimentation is part of the act of invention and until it has been undertaken, Applicant's claimed invention is incomplete. The instant situation is directly analogous to that which was addressed in *Brenner v. Manson*, 148 U.S.P.Q. 689 (Sus. Ct, 1966), in which the court held that:

Art Unit: 1646

“The basic quid pro quo contemplated by the Constitution and the Congress for granting a patent monopoly is the benefit derived by the public from an invention with substantial utility”, “[u]nless and until a process is refined and developed to this point-where specific benefit exists in currently available form-there is insufficient justification for permitting an applicant to engross what may prove to be a broad field”, and “a patent is not a hunting license”, “[i]t is not a reward for the search, but compensation for its successful conclusion.”

Applicant refers to three additional articles (Orntoft et al., Hyman et al. and Pollack et al.) as providing evidence that gene amplification generally results in elevated levels of the encoded polypeptide. Applicant characterizes Orntoft et al. as teaching in general (18 of 23 cases) chromosomal areas with more than 2-fold gain of DNA showed a corresponding increase in mRNA transcripts. Applicant characterizes Hyman et al. as providing evidence of a prominent global influence of copy number changes on gene expression levels. Applicant characterizes Pollack et al. as teaching that 62% of highly amplified genes show moderately or highly elevated expression and that, on average, a 2-fold change in DNA copy number is associated with a 1.5-fold change in mRNA levels. This has been fully considered but is not found to be persuasive. Orntoft et al. appear to have looked at increased DNA content over large regions of chromosomes and comparing that to mRNA and polypeptide levels from the chromosomal region. Their approach to investigating gene copy number was termed CGH. Orntoft et al. do not appear to look at gene amplification, mRNA levels and polypeptide levels from a single gene at a time. The instant specification reports data regarding amplification of individual genes, which are not likely to be in a chromosomal region which is highly amplified, given the low ΔCT values. Orntoft et al. concentrated on regions of chromosomes with strong gains of chromosomal material containing clusters of genes (p. 40). This analysis was not done for PRO274 in the instant specification. That is, it is not clear whether or not PRO274 is in a gene

Art Unit: 1646

cluster in a region of a chromosome that is highly amplified. Therefore, the relevance of Orntoft et al. is not clear. Hyman et al. used the same CGH approach in their research. Less than half (44%) of highly amplified genes showed mRNA overexpression (abstract). Polypeptide levels were not investigated. Therefore, Hyman et al. also do not support utility of the claimed antibodies to the polypeptides. Pollack et al. also used CGH technology, concentrating on large chromosome regions showing high amplification (p. 12965). Pollack et al. did not investigate polypeptide levels. Therefore, Pollack et al. also do not support the asserted utility of the claimed invention. Importantly, none of the three papers reported that the research was relevant to identifying probes that can be used as cancer diagnostics. The three papers state that the research was relevant to the development of **potential** cancer therapeutics, but also clearly imply that much further research was needed before such therapeutics were in readily available form. Accordingly, the specification's assertions that the claimed PRO274 antibodies have utility in the fields of cancer diagnostics and cancer therapeutics are not substantial.

The Polakis declaration under 37 CFR 1.132 filed Sept. 14, 2004 is insufficient to overcome the rejection of claims 58-62 based upon 35 U.S.C. §§ 101 and 112, first paragraph, as set forth in the last Office action for the following reasons:

Applicant presents a declaration by Dr. Polakis filed with the response under 37 CFR 1.132. In the declaration, Dr. Polakis states that the primary focus of the Tumor Antigen Project was to identify tumor cell markers useful as targets for cancer diagnostics and therapeutics. Dr. Polakis states that approximately 200 gene transcripts were identified that are present in human tumor cells at significantly higher levels than in corresponding normal human cells. Dr. Polakis states that antibodies to approximately 30 of the tumor antigen polypeptides have been

Art Unit: 1646

developed and used to show that approximately 80% of the samples show correlation between increased mRNA levels and changes in polypeptide levels. Dr. Polakis states that it remains a central dogma in molecular biology that increased mRNA levels are predictive of corresponding increased levels of the encoded polypeptide. Dr. Polakis characterizes the reports of instances where such a correlation does not exist as exceptions to the rule. This has been fully considered but is not found to be persuasive. First, it is important to note that the instant specification provides no information regarding increased mRNA levels of PRO274 in tumor samples relevant to normal samples. Only gene amplification data was presented. Therefore, the declaration is insufficient to overcome the rejection of claims 58-62 based upon 35 U.S.C. §§ 101 and 112, first paragraph, since it is limited to a discussion of data regarding the correlation of mRNA levels and polypeptide levels, and not gene amplification levels and polypeptide levels. Furthermore, the declaration does not provide data such that the examiner can independently draw conclusions. Only Dr. Polakis' conclusions are provided in the declaration. There is no evidentiary support to Dr. Polakis' statement that it remains a central dogma in molecular biology that increased mRNA levels are predictive of corresponding increased levels of the encoded polypeptide. Finally, it is noted that the literature cautions researchers from drawing conclusions based on small changes in transcript expression levels between normal and cancerous tissue. For example, Hu et al. (2003, Journal of Proteome Research 2:405-412) analyzed 2286 genes that showed a greater than 1-fold difference in mean expression level between breast cancer samples and normal samples in a microarray (p. 408, middle of right column). Hu et al. discovered that, for genes displaying a 5-fold change or less in tumors compared to normal, there was no evidence of a correlation between altered gene expression and

Art Unit: 1646

a known role in the disease. However, among genes with a 10-fold or more change in expression level, there was a strong and significant correlation between expression level and a published role in the disease (see discussion section). PRO 274 does not display a 10-fold or greater amplification, according to the specification.

Applicants further assert that even if one assumes that it is more likely than not that there is no correlation between gene amplification and increased mRNA/protein expression, a polypeptide encoded by a gene that is amplified in cancer would still have a specific and substantial utility, and provides the declaration by Dr. Avi Ashkenazi. Dr. Ashkenazi explains that even when amplification of a cancer marker gene does not result in significant over-expression of the corresponding gene product, this very absence of gene product over-expression still provides significant information for cancer diagnosis and treatment, in that if the gene product is over-expressed in some tumor types but not others, this would enable more accurate tumor classification and hence better determination of suitable therapy, and additionally, if a gene is amplified by the corresponding gene product is not-overexpressed, the clinician accordingly will decide not to treat a patient with agents that target that gene product

The declaration filed under 37 CFR 1.132 filed Sept. 14, 2004 is insufficient to overcome the rejection of claims 58-62 based upon lack of utility as set forth in the last Office action because: it has not been demonstrated that the protein of the instant invention is differentially expressed in different tumors. If it was, the protein would have a specific and substantial utility for tumor classification, but the mere assertion that it may be differentially expressed does not provide a specific and substantial utility, and is an invitation to experiment. The argument that if a gene is amplified but the gene product is not over-expressed, the clinician would accordingly

Art Unit: 1646

will decide not to treat a patient with agents that target the gene product is also insufficient to overcome the rejection of the claims. If a specific gene product was known to be involved in cancer and if there were known compounds that could be used to target the gene product, this would be an acceptable utility. However, the gene product of the instant invention has not been demonstrated to be involved in cancer. Over-expression of a gene product in a cancer cell does not necessarily mean that the gene product is involved in the cancer and that targeting the gene product would be therapeutic. Additionally, there are no known compounds that would target the gene product.

Applicants provide the Hanna et al. reference to support the Declaration of Dr. Ashkenazi. The Hanna reference is not applicable to the instant fact situation, as it deals with a known tumor associated gene, and not with a prospective analysis of the type found in this specification.

The proposed uses of the claimed invention are simply starting points for further research and investigation into potential practical uses of the claimed polypeptides. For all of these reasons, the rejections are maintained.

The following is a quotation of the first paragraph of 35 U.S.C. 112:

The specification shall contain a written description of the invention, and of the manner and process of making and using it, in such full, clear, concise, and exact terms as to enable any person skilled in the art to which it pertains, or with which it is most nearly connected, to make and use the same and shall set forth the best mode contemplated by the inventor of carrying out his invention.

4. Claims 58-62 also remain rejected under 35 U.S.C. 112, first paragraph. Specifically, since the claimed invention is not supported by either a specific and substantial asserted utility or a well established utility for the reasons set forth above, one skilled in the art clearly would not know how to use the claimed invention.

Rejections over Prior Art

Claim Rejections - 35 USC § 102 and § 103

The text of those sections of Title 35, U.S. Code not included in this action can be found in a prior Office action.

5.1 Claims 58-62 remain rejected under 35 U.S.C. 102(b) as being anticipated by Ho et al., Science, Vol. 289, July 14, 2000, pages 265-270, for reasons of record in the previous Office Action, mailed May 20, 2004, at page 10, and below.

5.2 Claims 59-62 remain rejected under 35 U.S.C. 103(a) as being unpatentable over Ho et al., Science, Vol. 289, July 14, 2000, pages 265-270, in view of Immunobiology, The Immune System in Health and Disease, Third Edition, Janeway, And Travers, Ed., 1997, for reasons of record in the previous Office Action, mailed May 20, 2004, at page 10, and below.

Applicants traverse rejections and assert that they rely on the gene amplification assay for patentable utility which was first disclosed in International Application no. PCT/US00/03565, filed Feb. 11, 2000, and assert that they are entitled to at least that filing date, so that Ho et al. is not prior art. Applicants' arguments have been fully considered but are not deemed persuasive, because the gene amplification assay fails to provide a patentable utility for the antibodies to the protein, for reasons discussed above.

It is believed that all pertinent arguments have been answered.

Conclusion

6. No claim is allowed.

Art Unit: 1646

THIS ACTION IS MADE FINAL. Applicant is reminded of the extension of time policy as set forth in 37 CFR 1.136(a).

A shortened statutory period for reply to this final action is set to expire **THREE MONTHS** from the mailing date of this action. In the event a first reply is filed within **TWO MONTHS** of the mailing date of this final action and the advisory action is not mailed until after the end of the **THREE-MONTH** shortened statutory period, then the shortened statutory period will expire on the date the advisory action is mailed, and any extension fee pursuant to 37 CFR 1.136(a) will be calculated from the mailing date of the advisory action. In no event, however, will the statutory period for reply expire later than **SIX MONTHS** from the mailing date of this final action.

Any inquiry concerning this communication or earlier communications from the examiner should be directed to Eileen B. O'Hara, whose telephone number is (571) 272-0878. The examiner can normally be reached on Monday through Friday from 10:00 AM to 6:30 PM.

If attempts to reach the examiner by telephone are unsuccessful, the examiner's supervisor, Brenda Brumback can be reached at (571) 272-0961.

The fax phone number for the organization where this application or proceeding is assigned is 703-872-9306.

Any inquiry of a general nature or relating to the status of this application should be directed to the Group receptionist whose telephone number is (571) 272-1600.

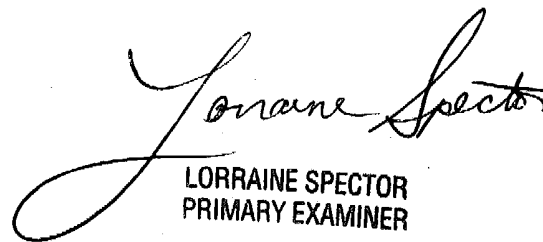
Information regarding the status of an application may be obtained from the Patent Application Information Retrieval (PAIR) system. Status information for published applications

Art Unit: 1646

may be obtained from either Private PAIR or Public PAIR. Status information for unpublished applications is available through Private PAIR only. For more information about the PAIR system, see <http://portal.uspto.gov/external/portal/pair>. Should you have questions on access to the Private PAIR system, contact the Electronic Business Center (EBC) at 866-217-9197 (toll free).

Eileen B. O'Hara, Ph.D.

Patent Examiner



LORRAINE SPECTOR
PRIMARY EXAMINER

Notice of References Cited	Application/Control No. 09/978,192	Applicant(s)/Patent Under Reexamination ASHKENAZI ET AL.	
	Examiner Eileen O'Hara	Art Unit 1646	Page 1 of 1

U.S. PATENT DOCUMENTS

*		Document Number Country Code-Number-Kind Code	Date MM-YYYY	Name	Classification
	A	US-			
	B	US-			
	C	US-			
	D	US-			
	E	US-			
	F	US-			
	G	US-			
	H	US-			
	I	US-			
	J	US-			
	K	US-			
	L	US-			
	M	US-			

FOREIGN PATENT DOCUMENTS

*		Document Number Country Code-Number-Kind Code	Date MM-YYYY	Country	Name	Classification
	N					
	O					
	P					
	Q					
	R					
	S					
	T					

NON-PATENT DOCUMENTS

*		Include as applicable: Author, Title Date, Publisher, Edition or Volume, Pertinent Pages)
	U	, Hu et al. (2003, Journal of Proteome Research 2:405-412.
	V	
	W	
	X	

*A copy of this reference is not being furnished with this Office action. (See MPEP § 707.05(a).)
Dates in MM-YYYY format are publication dates. Classifications may be US or foreign.

Analysis of Genomic and Proteomic Data Using Advanced Literature Mining

Yanhui Hu, Lisa M. Hines, Haifeng Weng, Dongmei Zuo, Miguel Rivera,
Andrea Richardson, and Joshua LaBaer*

Institute of Proteomics, Harvard Medical School—BCMP, 240 Longwood Avenue, Boston, Massachusetts 02115

Received March 13, 2003

High-throughput technologies, such as proteomic screening and DNA micro-arrays, produce vast amounts of data requiring comprehensive analytical methods to decipher the biologically relevant results. One approach would be to manually search the biomedical literature; however, this would be an arduous task. We developed an automated literature-mining tool, termed MedGene, which comprehensively summarizes and estimates the relative strengths of all human gene–disease relationships in Medline. Using MedGene, we analyzed a novel micro-array expression dataset comparing breast cancer and normal breast tissue in the context of existing knowledge. We found no correlation between the strength of the literature association and the magnitude of the difference in expression level when considering changes as high as 5-fold; however, a significant correlation was observed ($r = 0.41$; $p = 0.05$) among genes showing an expression difference of 10-fold or more. Interestingly, this only held true for estrogen receptor (ER) positive tumors, not ER negative. MedGene identified a set of relatively understudied, yet highly expressed genes in ER negative tumors worthy of further examination.

Keywords: bioinformatics • micro-array • text mining • gene-disease association • breast cancer

Introduction

At its current pace, the accumulation of biomedical literature outpaces the ability of most researchers and clinicians to stay abreast of their own immediate fields, let alone cover a broader range of topics. For example, to follow a single disease, e.g., breast cancer, a researcher would have had to scan 130 different journals and read 27 papers per day in 1999.¹ This problem is accentuated with high-throughput technologies such as DNA micro-arrays and proteomics, which require the analysis of large datasets involving thousands of genes, many of which are unfamiliar to a particular researcher. In any microarray experiment, thousands of genes may demonstrate statistically significant expression changes, but only a fraction of these may be relevant to the study. The ability to interpret these datasets would be enhanced if they could be compared to a comprehensive summary of what is known about all genes. Thus, there is a need to summarize existing knowledge in a format that allows for the rapid analysis of associations between genes and diseases or other specific biological concepts.

One solution to this problem is to compile structured digital resources, such as the Breast Cancer Gene Database¹ and the Tumor Gene Database.² However, as these resources are hand-curated, the labor-intensive review process becomes a rate-limiting step in the growth of the database. As a result, these

databases have a limited scale and the genes are not selected in a systematic fashion.

An alternative approach is automated text mining; a method which involves automated information extraction by searching documents for text strings and analyzing their frequency and context. This approach has been used successfully in several instances for biological applications. In most cases, it has been applied to extract information about the relationships or interactions that proteins or genes have with one another, in the literature or by functional annotation.^{3–7} Thus far, few publications have applied text-mining to examine the global relationships between genes and diseases. Perez-Iratxeta et al. automatically examined the GO (Gene Ontology) annotation of genes and their predicted chromosomal locations in order to identify genes linked to inherited disorders.⁸

To obtain a more global understanding of disease development, it would be valuable to incorporate information regarding all possible gene-disease relationships, including biochemical, physiological, pharmacological, epidemiological, as well as genetic. This information would enable comprehensive comparisons between large experimental datasets and existing knowledge in the literature. This would accomplish two things. First, it would serve to validate experiments by demonstrating that known responses occur as predicted. Second, it would rapidly highlight which genes are corroborated by the literature and which genes are novel in a given context. We have utilized a computational approach to literature mining to produce a

* To whom correspondence should be addressed: jlabar@hms.harvard.edu.

comprehensive set of gene-disease relationships. In addition, we have developed a novel approach to assess the strength of each association based on the frequency of citation and co-citation. We applied this tool to help interpret the data from a large micro-array gene expression experiment comparing normal and cancerous breast tissue.

Methods

MedGene Database. MedGene is a relational database, storing disease and gene information from NCBI, text mining results, statistical scores, and hyperlinks to the primary literature. MedGene has a web-based user interface for users to query the database (<http://hipseq.med.harvard.edu/MedGene/>).

Text Mining Algorithms. MeSH files were downloaded from the MeSH web site at NLM (National Library of Medicine) (<http://www.nlm.nih.gov/mesh/meshhome.html>) and human disease categories were selected. LocusLink files were downloaded from the LocusLink web site at NCBI (<http://www.ncbi.nih.gov/LocusLink/>). Official/preferred gene symbol, official/preferred gene name, and gene alternative symbols and names, all relevant annotations and URLs for each LocusLink record, were collected. Gene search terms were used for literature searching and included all qualified gene names, gene symbols, and gene family terms. Primary gene keys, predominantly qualified gene family terms and gene official/preferred symbols, were used to index Medline records. If the official/preferred gene symbols did not meet the standards to be an index, then qualified gene official/preferred names were used. A local copy of Medline records (up to July, 2002) was pre-selected.

A JAVA module examined the MeSH terms and then indexed each Medline record with the appropriate disease terms. A separate JAVA module was used to examine the titles and abstracts for gene search terms and then to index the gene-related Medline records with the relevant primary gene key(s).

Statistical Methods. For every gene and disease pair, we counted records that were indexed for both gene and disease (double positive hits), for disease only (disease single hits), for gene only (gene single hits), and for neither gene nor disease (double negative hits) to generate a 2×2 contingency table. On the basis of the contingency table-framework, we applied different statistical methods to estimate the strength of gene-disease relationships and evaluated the results. These methods included chi-square analysis, Fisher's exact probabilities, relative risk of gene, and relative risk of disease¹⁶ (<http://hipseq.med.harvard.edu/MedGene/>). In addition, we computed the "product of frequency", which is the product of the proportion of disease/gene double hits to disease single hits and the proportion of disease/gene double hits to gene single hits. To obtain a normal distribution, we transformed all the statistical scores using the natural logarithm. We selected the log of the product of frequency (LPF) to validate MedGene and to use for the analysis with the micro-array data. Spearman rank-correlation coefficients were used to assess the linear relationship between LPF and micro-array fold change in expression level.

Global Analysis. Diseases with at least 50 related genes were selected for clustering analysis, and the LPF scores were normalized with total score for each disease. Hierarchical clustering was done with the "Cluster" software and the clustering result was visualized using "TreeView" (<http://rana.lbl.gov/EisenSoftware.htm>).

Breast Tissue Micro-Arrays. Eighty-nine breast cancer samples (79% ER-positive) and 7 normal breast tissue samples were selected from the Harvard Breast SPORE frozen tissue repository and were representative of the spectrum of histological types, grades, and hormone receptor immuno-phenotypes of breast cancer. Biotinylated cRNA, generated from the total RNA extracted from the bulk tumor, was hybridized to Affymetrix U95A oligo-nucleotide micro-arrays. These micro-arrays consist of 12 400 probes, which represent approximately 9000 genes. Raw expression values were obtained using GENE-CHIP software from Affymetrix, and then further analyzed using the DNA-Chip Analyzer (dChip) custom software.

Results

Automated Indexing of Medline Records by Disease and Gene. To study the gene-disease associations in the literature, we first compiled complete lists for human diseases and human genes. To index all Medline records that were relevant to human diseases, the Medical Subject Heading (MeSH) index of Medline records was utilized. MeSH is a controlled medical vocabulary from the National Library of Medicine and consists of a set of terms or subject headings that are arranged in both an alphabetic and an hierarchical structure. Medline records are reviewed manually and MeSH terms are added to each with software assistance.^{9,10} Twenty-three human disease category headings along with all of their child terms (see the Supporting Information, Supplemental Table 1, or visit http://hipseq.med.harvard.edu/MedGene/publication/s_Table1.html) were selected from the 2002 MeSH Index creating a list of 4033 human diseases.

No index comparable to the MeSH index exists for genes, and thus, it was necessary to apply a string search algorithm for gene names or symbols found in Medline text. A complete list of genes, gene names, gene symbols, and frequently used synonyms were collected from the LocusLink database at NCBI,^{11,12} which contains 53 259 independent records keyed by an official gene symbol or name (June 18th, 2002). For the purposes of this study, no distinction was made between genes and their gene products. Authors often use the same name for both, differentiating the two only by the use of italics, if at all. For the intended use of this study, this lack of distinction is unlikely to have a large effect and may in fact be beneficial.

Initial attempts to search the literature using these lists revealed several sources of false positives and false negatives (Table 1). False positives primarily arose when the searched term had other meanings, whereas false negatives arose from syntax discrepancies necessitating the development of filters to reduce these errors. The syntax issues were readily handled by including alternate syntax forms in the search terms. The false positive cases, caused by duplicative and unrelated meanings for the terms, were more difficult to manage. Where possible, case sensitive string mapping reduced inappropriate citations. In many cases, however, this was not sufficient and the terms had to be eliminated entirely, thereby reducing the false positive rate but unavoidably under-representing some genes.

For the purposes of data tracking, a primary gene key was selected to represent all synonyms that correspond to each gene. Medline records were indexed with a primary gene key when any synonym for that key was found in the title or abstract. Case-insensitive string mapping was used for all searches except as noted above. No additional weight was

Table 1. Systematic Sources of False Positives and False Negatives in Unfiltered Data^a

source of error	error type	example	filter solution
gene symbol/name is not unique	false positive	MAG—myelin associated glycoprotein MAG—malignancy-associated protein	eliminate this term
gene symbol is unrelated abbreviation	false positive	PA—pallid homologue (mouse), pallidin (also abbrev. for Pennsylvania)	eliminate this term
gene symbol/name has language meaning	false positive	WAS—Wiskott–Aldrich Syndrome (also the word “was”)	case-sensitive string search
nonstandard syntax	false negative	BAG-1 instead of BAG1	add dash term
unofficial gene name/symbol	false negative	P53 instead of TP53	add all gene nicknames
nonspecified gene name	false negative	estrogen receptor instead of Estrogen receptor 1	add family stem term

^a In preliminary studies, Medline was searched for co-occurrence of genes and diseases and the resulting output was evaluated to identify error sources that were amenable to global filters. Each error source is categorized by the type of error it causes: false positives are suggested relationships that are not real and false negatives are real relationships that are underrepresented. The filter solutions used are indicated. Note that in some cases, the filter solution itself introduces error. In general, error rates maximized sensitivity, even at the expense of specificity if needed.

added for multiple occurrences of a term or the co-occurrence of multiple synonyms for the same gene key.

Medline records were searched with all qualified gene identifiers, such as the official/preferred gene symbol, the official/preferred gene name, all gene nicknames and all syntax variants. In situations where there are several members of a gene family or splice variants, some authors prefer to use a shortened gene family name, e.g., estrogen receptor instead of estrogen receptor 1 (*ESR1*), creating a source of false negatives. For this reason, gene family stem terms were created for all genes that have an alpha or numerical suffix (e.g., *IL2RA*, *TGFβ*, *ESR1*, etc.) and then used to search the literature. The family stem terms were handled separately from the specific gene names so that it would be clear when linkages were made to the gene family versus a specific member in that family.

To improve performance and accuracy, some pre-selection was applied to the records that were scanned. First, review articles were eliminated to avoid redundant treatment of citations. Second, non-English journals were removed because the natural language filters were only relevant to English publications. Finally, journals unlikely to contain primary data about gene-disease relationships were also removed (e.g., *Int. J. Health Educ.*, *Bedside Nurse*, and *J. Health Econ.*). Together, these filters reduced the 12 198 221 Medline publications (July 2002) by 37%.

Ranking the Relative Strengths of Gene-Disease Associations. In total, there were 618 708 gene-disease co-citations, in which 16% (8297) of all studied genes had been associated to a disease and 96% (3875) of all diseases had been associated to at least one gene. To rank the relative strengths of gene disease relationships, we tested several different statistical methods and examined the results. With the exception of the relative risk estimates, the methods provided similar results with respect to the rank order of the gene-disease association strengths. However, after comparing the results to other databases and after consulting disease experts, the log of the product of frequency (LPF) was selected for further analysis because it gave the best results overall.

Validation of MedGene. In developing this tool, it was important to minimize the number of missed genes (false negatives) and misclassified genes (false positives). However, in situations when these goals were in conflict, inclusiveness was prioritized. To determine the false negative rate in MedGene, breast cancer was used as a test case because it was associated with more genes than any other human disease and because

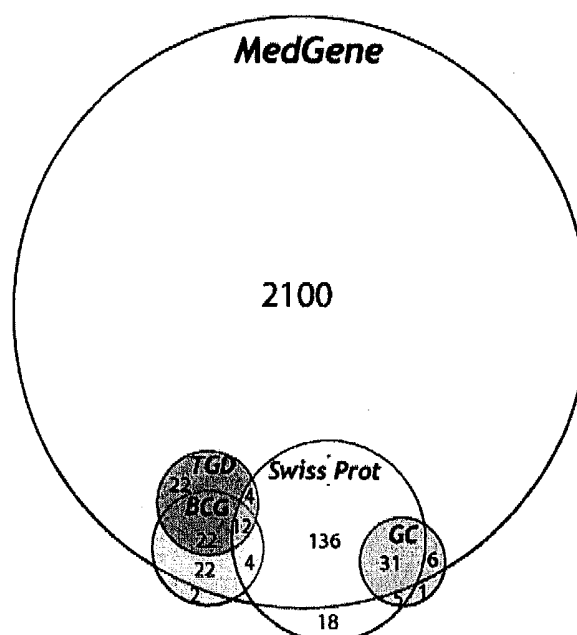


Figure 1. Estimation of the false negative rate by comparison with hand-curated databases. The breast cancer-related genes identified by MedGene were compared with those listed in several other databases including the Tumor Gene Database (TGD),² the Breast Cancer Gene Database (BCG),¹ GeneCards (GC)¹⁷ and Swissprot.¹⁸ Genes were considered false negatives if they were represented in at least one of these other databases and not in MedGene and their link to breast cancer was supported by at least one literature reference. All literature references were verified by manual review to confirm their validity. The number of genes in each database or shared by more than one database is indicated. The false negative rate was calculated by genes missed at MedGene (26)/total number of nonoverlapping genes in other databases (285).

there were several public databases that link genes to breast cancer. We compared the list of breast cancer-related genes from MedGene to these databases, illustrated in Figure 1. Among the 285 distinct breast cancer-related genes that were supported by at least one literature citation in these hand-curated databases, 26 were absent from MedGene, suggesting a false negative rate of approximately 9%. To determine why these were missed, all literature references for these genes (80

papers) were reviewed manually (see the Supporting Information, Supplemental Table 2, or visit http://hipseq.med.harvard.edu/MedGene/publication/s_Table_2.html). Among these papers, most false negatives were caused by nonstandard gene terms or gene terms eliminated by our specificity filters. Few genes were missed because they were only mentioned in review papers (0.4%) or they appeared only in the body of the manuscript but not the abstract or title (1.1%). Of note, MedGene identified approximately 2000 additional breast cancer-related genes not listed in any other database.

To assess the false positive error rate, two complementary approaches were used: a detailed analysis of one disease and a global examination of 1000 diseases. The detailed approach examined the false positive error rate and its sources, whereas the global approach tested whether the overall results made biomedical sense.

Using the LPF, 1467 genes related to prostate cancer were assembled in rank order. We then retrieved approximately 300 Medline records each for the highest ranked 100 and the lowest ranked 200 genes and manually reviewed the titles and abstracts to determine the verity of the association. Nearly 80% of the highest ranked 100 genes fell into one of the five categories that reflect meaningful gene-disease relationships (see the Supporting Information, Supplemental Table 3, or visit http://hipseq.med.harvard.edu/MedGene/publication/s_Table_3.html). Among the lowest ranked 200 genes, approximately 70% reflected true relationships. Of the 600 records reviewed, there were only two in which the association between the gene and the disease was described as negative. Both were genes with very low scores. In both cases, the authors did not argue the absence of any relationship, but rather that a particular feature of the gene or protein was not shown to be related to human prostate cancer.^{13,14}

The coincidence of some gene symbols with medical abbreviations, chemical abbreviations and biological abbreviations resulted in most of the false positives (see the Supporting Information, Supplemental Table 4, or visit http://hipseq.med.harvard.edu/MedGene/publication/s_Table_4.html), emphasizing the importance of the filters that were added in the search algorithm (Table 1). Without the filters, the false positive rate more than doubled, and the false negative rate rose dramatically (data not shown). For example, among the papers about breast cancer, there were only 12 Medline records that referred to *ESR1* and 10 to *ESR2*, whereas almost 2000 papers mentioned estrogen receptor without specifying *ESR1* or *ESR2*; this latter group was detected by the family stem term filter.

To further validate these results, a global analysis of the gene-disease relationships described by MedGene was performed. For this experiment, it was reasoned that the more closely related the diseases are to one another, the more they will be related to the same gene sets. Thus, if the relationships defined by MedGene accurately reflected the literature, then an unsupervised hierarchical clustering of the gene data should group diseases in a manner consistent with common medical thinking. Conversely, if the clustered diseases do not make sense biologically or medically, it may reflect excessive false positives, false negatives, or inappropriate scoring of the data.

To execute this experiment, the gene sets and the corresponding LPF values for 1000 randomly selected diseases (each with at least 50 gene relationships) were used as a dataset for clustering the diseases. A review of the results showed that the resulting disease clusters were indeed logical based upon common medical knowledge (see the Supporting Information,

Supplemental Figure 1, or visit http://hipseq.med.harvard.edu/MedGene/publication/s_Figure_1.html). For example, in one such cluster shown in Figure 2, diabetes and its complications grouped together and were also closely linked to diseases associated with starvation states.

The number of genes associated with a given disease can be estimated by adjusting the MedGene number up by the false negative rate (~9%) and down by the false positive rate (~26% on average). Using this, the average disease has 103.7 ± 45.3 (mean \pm s.d.) genes associated with it, although the range is quite broad with 2359 genes related to breast cancer, 2122 genes related to lung cancer and no genes related to a number of diseases.

Applying MedGene to the Analysis of Large Datasets. Access to a comprehensive summary of the genes linked to human diseases provided an opportunity to analyze data obtained from a high-throughput experiment. We compared the MedGene breast cancer gene list to a gene expression data set generated from a micro-array analysis comparing breast cancer and normal breast tissue samples. Micro-array analysis identified 2286 genes that had greater than a 1-fold difference in mean expression level between breast cancer samples and normal breast samples. Using MedGene, we sorted the 2286 genes into four classes: 555 genes directly linked to breast cancer in the literature by gene term search (first-degree association by gene name); 328 genes directly linked by family term search (first-degree association by family term); 1021 genes linked to breast cancer only through other breast cancer genes (second-degree association); and 505 genes not previously associated with breast cancer. (See the Supporting Information, Supplemental Figure 2, or visit http://hipseq.med.harvard.edu/MedGene/publication/s_Figure_2.html.) Among the 505 previously unrelated genes, 467 were either newly identified genes or genes that had not previously been associated with any disease. Among the remaining 38 genes, 9 had been related to other cancers, specifically esophageal, colon, uterine, skin, and cervix.

To determine whether the genes highlighted by the micro-array analysis were more likely to have been previously linked to breast cancer in the literature, we created a two-dimensional plot of the fold change of expression level between breast cancer and normal tissue versus the literature score (LPF) (Figure 3A). There was a broad spread of expression changes among the genes directly linked to breast cancer ranging from less than 1-fold change (68%) to over 40-fold (0.3%). Notably, the majority of genes with greater than 10-fold expression changes were linked to breast cancer by first-degree association.

Among all 754 genes directly linked to breast cancer in the literature, there was no correlation between LPF and micro-array fold change ($r = 0.018$, p -value = 0.62). However, when we stratified the analysis based on the magnitude of the fold change, we observed an increasing trend in correlation (Figure 3B) suggesting that genes with a more substantial change in expression level were more likely to have a stronger association in the literature. For genes that had 10-fold change or more in expression level, the correlation increased to 0.41 (p -value = 0.05).

When we evaluated the micro-array data separately for ER positive and ER negative tumors, the trend in correlation between fold change and literature score was highly dependent on estrogen receptor status. Interestingly, there was a similar trend in correlation for ER positive tumors, but no trend in correlation for ER negative tumors.

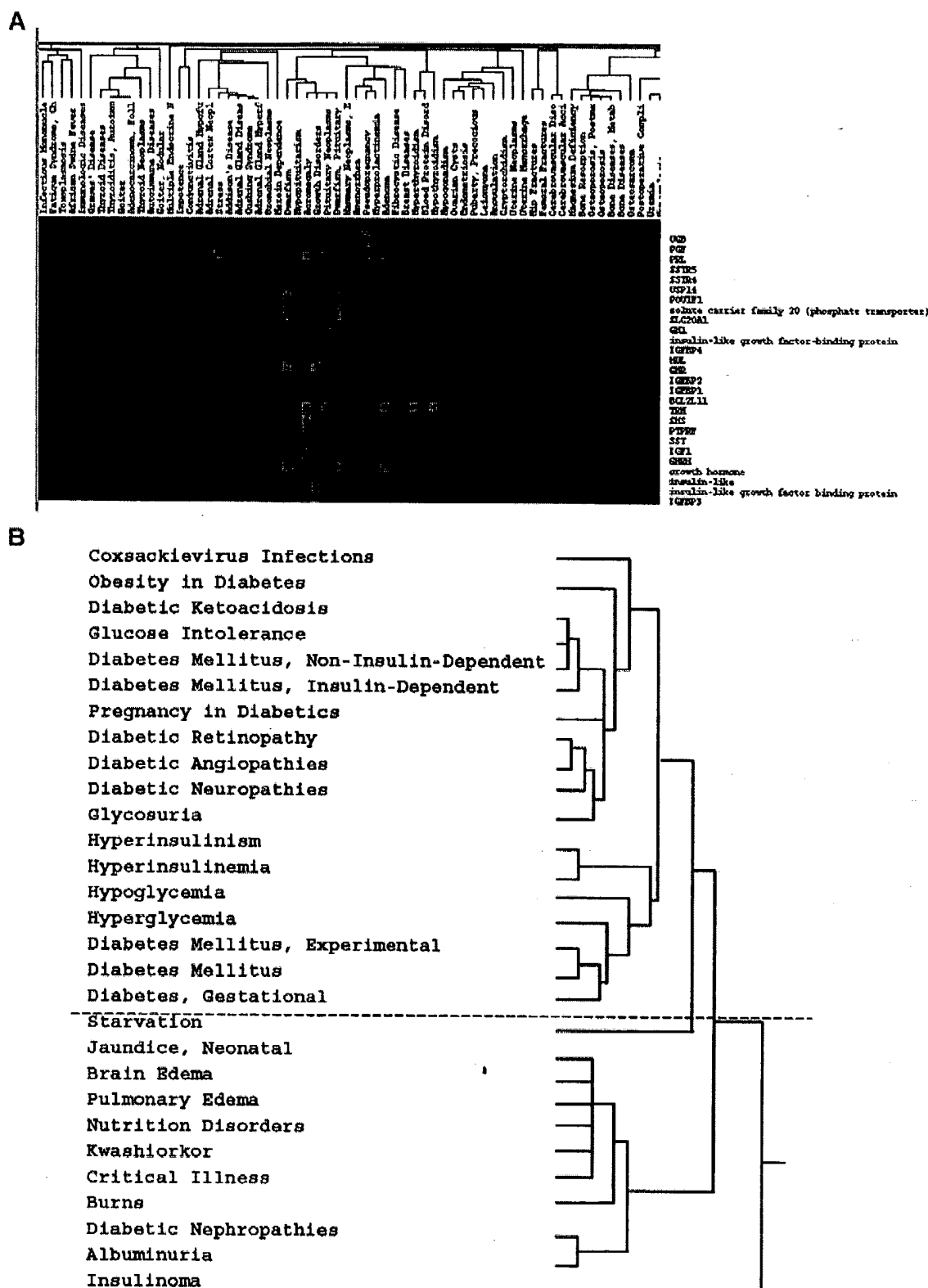


Figure 2. Global validation by clustering analysis. 2(A). The gene sets and the corresponding LPF values for 1000 diseases, each with at least 50 gene relationships, were used in an unsupervised clustering of the diseases based on the gene patterns associated with them. A sample of the data is shown here. 2(B). One of the resulting clusters is shown that corresponds to blood sugar states. Diabetes terms (above the line) and starvation states terms (under the line) clustered together. Within these groups, there is also clustering of diabetic small vessel complications, altered serum chemistries, nutritional disorders, etc. (Supplemental Figure 1: http://hipseq.med.harvard.edu/MedGene/publication/s_Figure 1.html).

Finally, to validate our findings, we computed similar correlations between the breast cancer expression data and LPF scores generated by MedGene for hypertension, a

disease unrelated to breast cancer. As expected, we did not observe an increasing trend in correlation for hypertension.

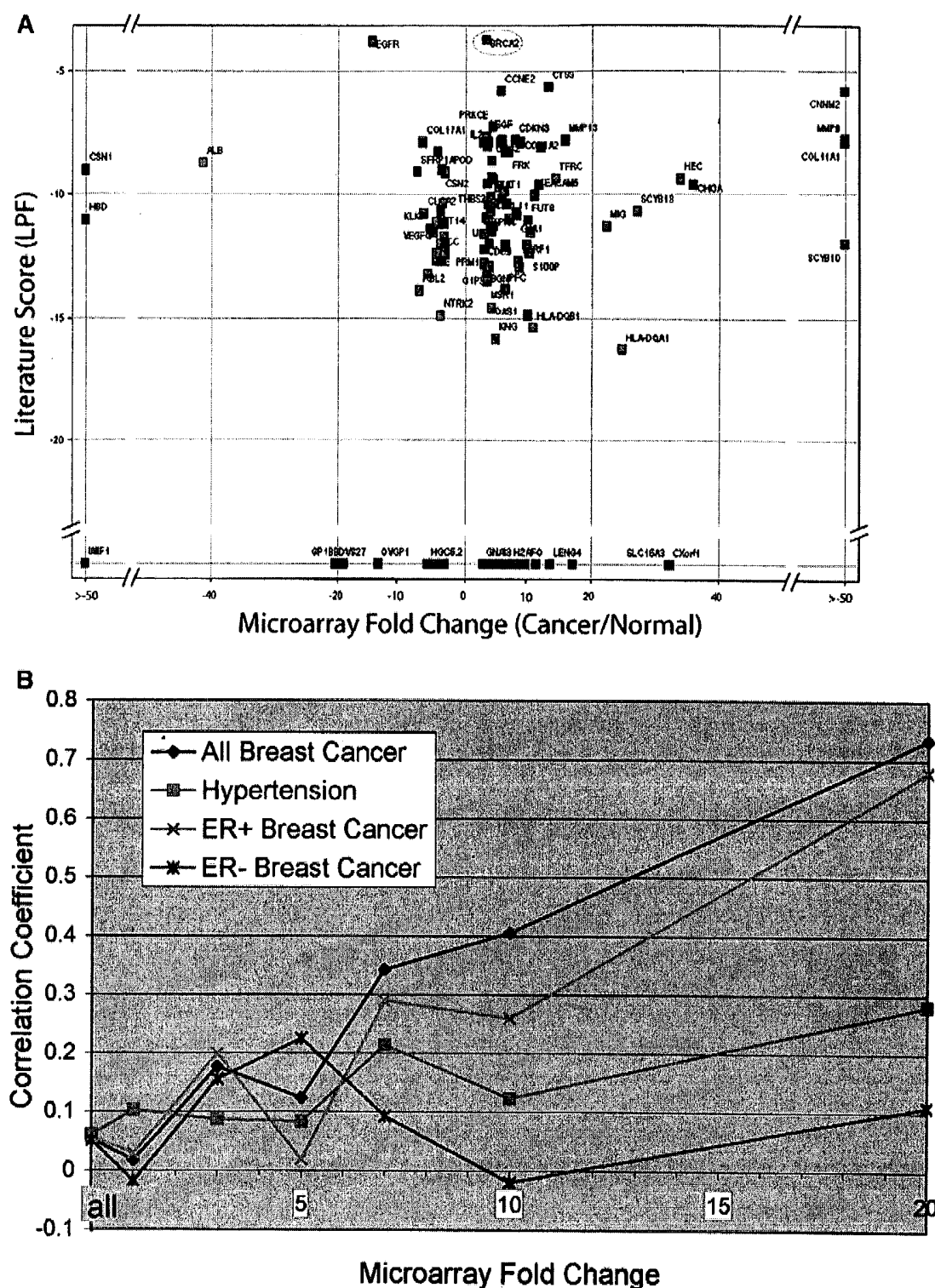


Figure 3. Relationship between literature score and functional data for breast cancer. **3A.** The data from an expression analysis of samples for breast tumors and normal breast tissue were analyzed to indicate the fold difference of expression level between breast tumor and normal sample (cutoff ≥ 3 -fold change). The fold changes were plotted against the literature score for the same gene set. Green dots represent first-degree association by gene search, blue dots represent first-degree association by family search and red dots represent no-association. Some well-studied genes, such as BRCA2 (pink circle), are not reflected by a substantial difference in expression level. Furthermore, the majority of genes that have no association with breast cancer in the literature had less than 10-fold expression changes (shaded area). **3B.** The Spearman rank-correlation coefficients between literature score (LPF) and the fold change of expression level between tumor and normal breast samples (y-axis) in relation to the amount of fold change of expression level (x-axis). Gene rank lists were generated for breast cancer (blue) and hypertension (pink). Correlations were also computed between the breast cancer gene LPF scores and fold change expression data among estrogen receptor positive tumors only (light blue) and estrogen receptor negative tumors only (purple).

Table 2. Top 25 Genes Related to Selected Human Diseases^a

breast neoplasms	hypertension	rheumatoid arthritis	bipolar disorder	atherosclerosis
estrogen receptor	<i>REN</i>	<i>RA</i>	<i>ERDA1</i>	apolipoprotein
<i>PCR</i>	<i>DBP</i>	<i>TNFRSF10A</i>	<i>SNAP29</i>	<i>APOE</i>
<i>ERBB2</i>	<i>LEP</i>	<i>CRP</i>	<i>PFKL</i>	<i>LDLR</i>
<i>BRCA1</i>	<i>AGT</i>	<i>AS</i>	<i>DRD2</i>	<i>ELN</i>
<i>BRCA2</i>	<i>INS</i>	<i>ESR1</i>	<i>TRH</i>	<i>ARG1</i>
<i>EGFR</i>	kallikrein	<i>HLA-DRB1</i>	<i>IMPA2</i>	<i>APOB</i>
<i>CYP19</i>	<i>ACE</i>	<i>DR1</i>	<i>HTR3A</i>	<i>APOA1</i>
<i>TFF1</i>	endothelin	interleukin	<i>DRD3</i>	<i>MSR1</i>
<i>PSEN2</i>	<i>S100A6</i>	<i>TNF</i>	<i>REM</i>	<i>LPL</i>
<i>TP53</i>	<i>BDK</i>	<i>IL6</i>	<i>KCNN3</i>	<i>PON1</i>
<i>CES3</i>	<i>DIANPH</i>	collagen	<i>DRD4</i>	plasminogen
<i>CEACAM5</i>	<i>SAR1</i>	<i>IL1A</i>	<i>HTR2C</i>	activator inhibitor
<i>ERBB3</i>	<i>PIH</i>	<i>ACR</i>	<i>RELN</i>	<i>PLG</i>
cyclin	<i>CD59</i>	<i>TNFRSF12</i>	<i>DBH</i>	vascular cell
<i>COX5A</i>	<i>ALB</i>	<i>IL2</i>	<i>MACA</i>	adhesion molecule
cathepsin	<i>CYP11B2</i>	<i>CHI3L1</i>	<i>COMT</i>	<i>ATOH1</i>
<i>ERBB4</i>	<i>MAT2B</i>	<i>IL8</i>	<i>HTR2A</i>	<i>VWF</i>
<i>TRAM</i>	angiotensin receptor	interleukin 1 matrix	<i>SYNJ1</i>	<i>INS</i>
<i>CCND1</i>	<i>AGTR2</i>	metalloproteinase	<i>INPP1</i>	<i>ARG2</i>
<i>EGF</i>	<i>NPPA</i>	interferon	<i>NEDD4L</i>	<i>ABCA1</i>
<i>MUC1</i>	<i>LVM</i>	<i>CD68</i>	<i>FRA13C</i>	<i>OLR1</i>
insulin-like	<i>DBH</i>	<i>IL4</i>	transducer of	collagen
<i>BCL2</i>	<i>NPY</i>	<i>IL17</i>	<i>ERBB2</i>	<i>MCP</i>
mucin	<i>POMC</i>	<i>MMP3</i>	<i>BAIAP3</i>	lipoprotein
<i>FGF3</i>	neuropeptide	<i>SIL</i>	<i>ATP1B3</i>	<i>APOA2</i>
			<i>DRD5</i>	intercellular
				adhesion molecule
				<i>RAB27A</i>

^a MedGene results for the top 25 genes associated with breast neoplasms, hypertension, rheumatoid arthritis, bipolar disorder, and atherosclerosis, respectively, ranked by LPF scores. The hyperlink to all the papers co-citing the gene and the disease is available at MedGene website (<http://hipseq.med.harvard.edu/MedGene/>).

Discussion

The Human Genome Project heralded a new era in biological research where the emphasis on understanding specific pathways has expanded to global studies of genomic organization and biological systems. High-throughput technologies can provide novel insight into comprehensive biological function but also introduces new challenges. The utility of these technologies is limited to the ability to generate, analyze, and interpret large gene lists. MedGene, a relational database derived by mining the information in Medline, was created to address this need. MedGene users can query for a rank-ordered list of human gene-disease relationships (Table 2) for one or more diseases. Each entry is hyperlinked to the original papers supporting each association and to other relevant databases.

MedGene is an innovative extension of previous text mining approaches. Perez-Iratxeta et al. used the GO annotation and their chromosomal locations to predict genes that may contribute to inherited disorders.⁸ MedGene takes a broader view and includes all diseases and all possible gene-disease relationships. Furthermore, MedGene utilizes co-citation to indicate a relationship rather than GO annotation, which is limited to the subset of genes that have GO annotation. Our approach is complementary to that taken by Chaussabel and Sher, who used the frequency of co-cited terms to cluster genes into a hierarchy of gene-gene relationships.⁶

A unique aspect of this tool is the ability to assess the relative strengths of gene-disease relationships based on the frequency of both co-citation and single citation. This presupposes that most co-citations describe a positive association, often referred to as publication bias¹⁵ and is supported by our observations

that negative associations are rare (Supplemental Table 3: http://hipseq.med.harvard.edu/MedGene/publication/s_Table3.html). Of course, relationships established by frequency of co-citation do not necessarily represent a true biological link; however, it is strong evidence to support a true relationship.

Another important feature of MedGene is the implementation of software filters that substantially reduced the error rate. We estimate that less than 10% of all associations were missed and at least 70% of even the weakest associations were real. For this study, all of the filters that we applied were general ones, e.g., expanding the list of all gene names to address the different syntax forms used by different journals, eliminating gene names that correspond to common English words, etc. The majority of the remaining search term ambiguities were idiosyncratic and difficult to identify systematically without causing a significant rise in false negatives. Alternative approaches, such as the examination of the nearest neighbor terms, need to be considered to further reduce the false positive rate.

It is not uncommon to see expression changes in microarray experiments as small as 2-fold reported in the literature. Even when these expression changes are statistically significant, it is not always clear if they are biologically meaningful. When comparing expression levels of disease to normal tissue, one expects an enrichment of known disease-related genes to appear in the altered expression group. MedGene provided a unique opportunity to test this notion in the context of existing knowledge on a novel breast cancer microarray dataset. For genes displaying a 5-fold change or less in tumors compared to normal, there was no evidence of a correlation between altered gene expression and a known role in the disease. This

Table 3. Genes with Large Expression Changes in ER- but Not in ER+ Breast Tumors

gene symbol	fold change (ER+)	fold change (ER-)
<i>KRTHB1</i>	1.0	610.8
<i>BRS3</i>	1.2	89.4
<i>DKK1</i>	1.2	69.8
<i>ZIC1</i>	1.9	59.6
<i>TLR1</i>	1.0	38.5
<i>KIAA0680</i>	2.6	33.2
<i>CDKN3</i>	1.0	30.6
<i>EBI2</i>	4.0	27.9
<i>GZMB</i>	3.8	21.9
<i>STK18</i>	4.7	18.6
<i>GPR49</i>	1.0	14.6
<i>MYO10</i>	1.6	14.4
<i>LAD1</i>	-1.0	13.5
<i>POLE2</i>	4.2	13.0
<i>HMG4</i>	4.4	12.9
<i>BCL2L11</i>	-1.2	12.3
<i>LRP8</i>	2.9	12.2
<i>CCNB2</i>	1.0	11.8
<i>CCNE2</i>	4.0	11.6
<i>FCB</i>	-4.3	11.1
<i>KNSL6</i>	2.9	10.9
<i>H1F5</i>	3.0	10.2
<i>SERPINH2</i>	4.6	10.2
<i>YAP1</i>	1.0	10.0
<i>IPHB</i>	-1.3	-10.4
<i>TCEA2</i>	-1.1	-10.8
<i>TFF1</i>	1.3	-11.4
<i>COL17A1</i>	-4.1	-15.7
<i>POP5</i>	1.1	-16.2
<i>BPAG1</i>	-4.6	-22.3
<i>PDZK1</i>	-1.1	-36.8
<i>VEGFC</i>	-2.8	-51.5
<i>MUC6</i>	-1.4	-64.9
<i>SERPINA5</i>	-1.0	-83.1
<i>MEIS1</i>	-1.6	-85.9
<i>CA12</i>	2.4	-150.3

Table 3. MedGene identified a set of relatively understudied, yet highly expressed genes in ER negative, but not ER positive breast tumors. All of these genes have either never been co-cited with breast cancer or have a weak association except those marked with an *.

reflects the many genes whose role in breast cancer may not involve large changes in expression in sporadic tumors (e.g., *BRCA1* and *BRCA2*) and genes whose modest changes in expression may be unrelated to the disease. Strikingly, among genes with a 10-fold change or more in expression level, there was a strong and significant correlation between expression level and a published role in the disease, providing the first global validation of the micro-array approach to identifying disease-specific genes.

The results derived from MedGene have two implications. First, a careful hunt for corroborating evidence of a role in breast cancer should precede any further study of genes with less than 5-fold expression level changes. Second, any genes with 10-fold changes or more are likely to be related to breast cancer and warrant attention. It is likely that this threshold will change depending on the disease as well as the experiment.

Interestingly, the observed correlation was only found among ER-positive tumors, not ER-negative. This may reflect a bias in the literature to study the more prevalent type of tumor in the population. Furthermore, this emphasizes that caution must be taken when interpreting experiments that may contain subpopulations that behave very differently. The MedGene approach identified a set of relatively understudied, yet highly expressed genes in ER-negative tumors that are worthy of further examination (Table 3).

In conclusion, we have developed an automated method of summarizing and organizing the vast biomedical literature. To our knowledge, the resulting database is the most comprehensive and accurate of its kind. By generating a score that reflects the strength of the association, it provides an important tool for the rapid and flexible analysis of large datasets from various high-throughput screening experiments. Furthermore, it can be used for selecting subsets of genes for functional studies, for building disease-specific arrays, for looking at genes common to multiple diseases and various other high-throughput applications. In the future, it will be possible to enhance the utility of the MedGene database by building links between genes and other MeSH terms as well as other biological processes and concepts, such as cell division and responses to small molecules.

Acknowledgment. We would like to thank P. Braun, L. Garraway, J. Pearlberg, and other members of our institute for helpful discussion. Many thanks to the NLM (National Library of Medicine) for licensing of MEDLINE and the annotation effort of adding MeSH indexes for MEDLINE abstracts. This work was funded by grants from the Breast Cancer Research Foundation and an NHLBI PGA Grant (Vol HL66582-02).

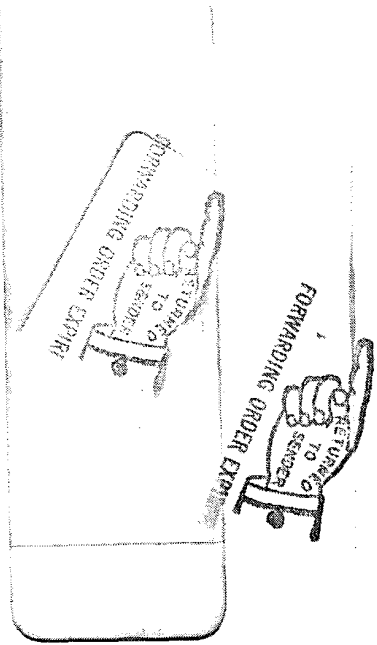
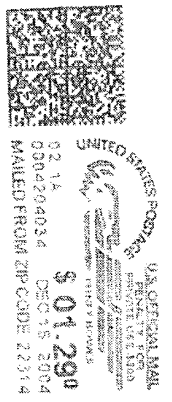
Supporting Information Available: Twenty-three human disease category headings along with all of their child terms selected from the 2002 MeSH index (Supplemental Table 1); analysis of the causes of false negatives in MedGene (Supplemental Table 2); meaningful gene-disease relationships found in MedGene (Supplemental Table 3); causes for incorrect assignment of gene indexes (Supplemental Table 4); a review of the results, showing that the resulting disease clusters were indeed logical (Supplemental Figure 1); and a review of the results showing that among the 505 previously unrelated genes, 467 were either newly identified genes or genes that had not previously been associated with any disease (Supplemental Figure 2). This material is available free of charge via the Internet at <http://pubs.acs.org> and at the web sites mentioned in the text.

References

- (1) Baasiri, R. A.; Glasser, S. R.; Steffen, D. L.; Wheeler, D. A. *Oncogene* **1999**, *18*, 7958–7965.
- (2) Steffen, D. L.; Levine, A. E.; Yarus, S.; Baasiri, R. A.; Wheeler, D. A. *Bioinformatics* **2000**, *16*, 639–649.
- (3) Marcotte, E. M.; Xenarios, I.; Eisenberg, D. *Bioinformatics* **2001**, *17*, 359–363.
- (4) Ono, T.; Hishigaki, H.; Tanigami, A.; Takagi, T. *Bioinformatics* **2001**, *17*, 155–161.
- (5) Jenssen, T. K.; Laegreid, A.; Komorowski, J.; Hovig, E. *Nat. Genet.* **2001**, *28*, 21–28.
- (6) Chaussabel, D.; Sher, A. *Genome Biol.* **2002**, *3*, RESEARCH0055.
- (7) Gibbons, F. D.; Roth, F. P. *Genome Res.* **2002**, *12*, 1574–1581.
- (8) Perez-Iratxeta, C.; Bork, P.; Andrade, M. A. *Nat. Genet.* **2002**, *31*, 316–319.
- (9) Funk, M. E.; Reid, C. A. *Bull. Med. Libr. Assoc.* **1983**, *71*, 176–183.
- (10) Humphrey, S. M.; Miller, N. E. *J. Am. Soc. Inf. Sci.* **1987**, *38*, 184–196.
- (11) Maglott, D. R.; Katz, K. S.; Sicotte, H.; Pruitt, K. D. *Nucleic Acids Res.* **2000**, *28*, 126–128.
- (12) Pruitt, K. D.; Maglott, D. R. *Nucleic Acids Res.* **2001**, *29*, 137–140.
- (13) Wadelius, M.; Andersson, A. O.; Johansson, J. E.; Wadelius, C.; Rane, E. *Pharmacogenetics* **1999**, *9*, 333–340.
- (14) Adam, R. M.; Borer, J. G.; Williams, J.; Eastham, J. A.; Loughlin, K. R.; Freeman, M. R. *Endocrinology* **1999**, *140*, 5866–5875.
- (15) Montori, V. M.; Smieja, M.; Guyatt, G. H. *Mayo Clin. Proc.* **2000**, *75*, 1284–1288.
- (16) Denenberg, V. H. *Statistics Experimental Design for Behavioral and Biological Researchers*; Wiley-Liss: New York, 1976.
- (17) Rebhan, M.; Chalfia-Caspi, V.; Prilusky, J.; Lancet, D. *Trends Genet.* **1997**, *13*, 163.
- (18) Balroch, A.; Apweiler, R. *Nucleic Acids Res.* **2000**, *28*, 45–48. PR0340227

Organization **IC1600** Bldg **Room 152N**
U. S. DEPARTMENT OF COMMERCE
COMMISSIONER FOR PATENTS
P.O. BOX 1450
ALEXANDRIA, VA 22313-1450
IF UNDELIVERABLE RETURN IN TEN DAYS
OFFICIAL BUSINESS

AN EQUAL OPPORTUNITY EMPLOYER



TECH CENTER 1600/2900

DEC 29 2004

RECEIVED